

SVD and Least Squares Problems

Note Title

LECTURE 17

★ LS via SVD

Recall the LS solution via QR factorization:

- (1) Compute reduced QR of A .
- (2) Compute $y = \hat{Q}^T b$.
- (3) Solve $\hat{R} x = y$ — (*)

If A : full rank, then $\hat{R}_{ii} \neq 0$, $1 \leq i \leq n$, and the triangular system (*) has a unique LS solution.

Now using the reduced SVD of A , i.e., $A = \hat{U} \hat{\Sigma} V^T$, we can also solve the normal eqn:

$$\begin{aligned} A^T A x &= A^T b \\ \Leftrightarrow (\hat{U} \hat{\Sigma} V^T)^T (\hat{U} \hat{\Sigma} V^T) x &= (\hat{U} \hat{\Sigma} V^T)^T b \\ \Leftrightarrow V \hat{\Sigma}^T \hat{U}^T \hat{U} \hat{\Sigma} V^T x &= V \hat{\Sigma} \hat{U}^T b \\ \Leftrightarrow V \hat{\Sigma}^T \hat{\Sigma} V^T x &= V \hat{\Sigma}^T \hat{U}^T b \\ \Leftrightarrow \hat{\Sigma}^T \hat{\Sigma} V^T x &= \hat{\Sigma}^T \hat{U}^T b \quad \text{since } V: \text{ortho.} \\ \Leftrightarrow \hat{\Sigma} V^T x &= \hat{U}^T b \quad \text{if } A: \text{full rank,} \\ &\quad \text{i.e., } \sigma_j > 0, 1 \leq j \leq n \end{aligned}$$

This can be solved easily.

- (1) Compute reduced SVD of A .
- (2) Compute $y = \hat{U}^T b$.
- (3) Solve $\hat{\Sigma} w = y$. — (**)
- (4) Set $x = V w$.

Note: (**) is a diagonal system, easier to solve than (*) !!

★ Pseudoinverse and SVD

Recall that if $A \in \mathbb{R}^{m \times n}$ is full rank,

$$\underline{m > n} : A^+ = (A^T A)^{-1} A^T$$

$$\underline{m = n} : A^+ = A^{-1}$$

$$\underline{m < n} : A^+ = A^T (A A^T)^{-1}$$

However, we can define the pseudo inv. using SVD even if A is not full rank!

$$A = U \Sigma V^T, \quad \Sigma = \begin{array}{c|c|c} \overbrace{\begin{matrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_r \end{matrix}}^r & \overbrace{\begin{matrix} 0 & & 0 \\ & \ddots & \\ 0 & & 0 \end{matrix}}^{n-r} & \\ \hline \begin{matrix} 0 & & 0 \end{matrix} & \begin{matrix} 0 & & 0 \end{matrix} & \\ \hline \end{array} \left. \begin{array}{l} \\ \\ \end{array} \right\} \begin{array}{l} r \\ m-r \end{array}$$

Define

$$A^+ := V \Sigma^+ U^T, \quad \Sigma^+ = \begin{array}{c|c|c} \overbrace{\begin{matrix} 1/\sigma_1 & & 0 \\ & \ddots & \\ 0 & & 1/\sigma_r \end{matrix}}^r & \overbrace{\begin{matrix} 0 & & 0 \\ & \ddots & \\ 0 & & 0 \end{matrix}}^{n-r} & \\ \hline \begin{matrix} 0 & & 0 \end{matrix} & \begin{matrix} 0 & & 0 \end{matrix} & \\ \hline \end{array} \left. \begin{array}{l} \\ \\ \end{array} \right\} \begin{array}{l} r \\ n-r \end{array}$$

As we discussed before, A^+ satisfies the following **Moore-Penrose conditions**:

$$(i) \quad A X A = A; \quad (ii) \quad X A X = X$$

$$(iii) \quad (A X)^T = A X; \quad (iv) \quad (X A)^T = X A.$$

Such X is uniquely determined and $X = A^+ !!$

★ Pseudoinverse & Orthogonal Projectors

Thm AA^T is an ortho. proj. onto $\text{range}(A)$

$$\text{and } AA^T = U_r U_r^T$$

$A^T A$ is an ortho. proj. onto $\text{range}(A^T)$

$$\text{and } A^T A = V_r V_r^T$$

where $U_r \in \mathbb{R}^{m \times r}$, $V_r \in \mathbb{R}^{n \times r}$ consist of the first r columns of U, V , respectively.
 $r = \text{rank}(A)$.

(Proof) Let $P_A := AA^T$, $P_{A^T} := A^T A$.

$$\text{Now, } P_A = U \Sigma V^T V \Sigma^+ U^T$$

$$= U \Sigma \Sigma^+ U^T = U \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix} U^T$$

$$= U_r U_r^T \quad \checkmark$$

$$P_A^2 = U_r \underbrace{U_r^T U_r}_{= I_r} U_r U_r^T = U_r U_r^T = P_A \quad \checkmark$$

so it's a proj.!

$$P_A^T = (U_r U_r^T)^T = (U_r^T)^T U_r = U_r U_r^T = P_A \quad \checkmark$$

so it's an ortho. proj.!

Finally, it's also clear that

P_A maps onto $\text{range}(A)$ since

$$\text{range}(A) = \langle u_1, \dots, u_r \rangle. \quad \checkmark$$

You can do similarly for P_{A^T} ///

Note: Consider any $x \in \text{range}(A)$.

Then $\exists y \in \mathbb{R}^n$ s.t. $x = Ay$.

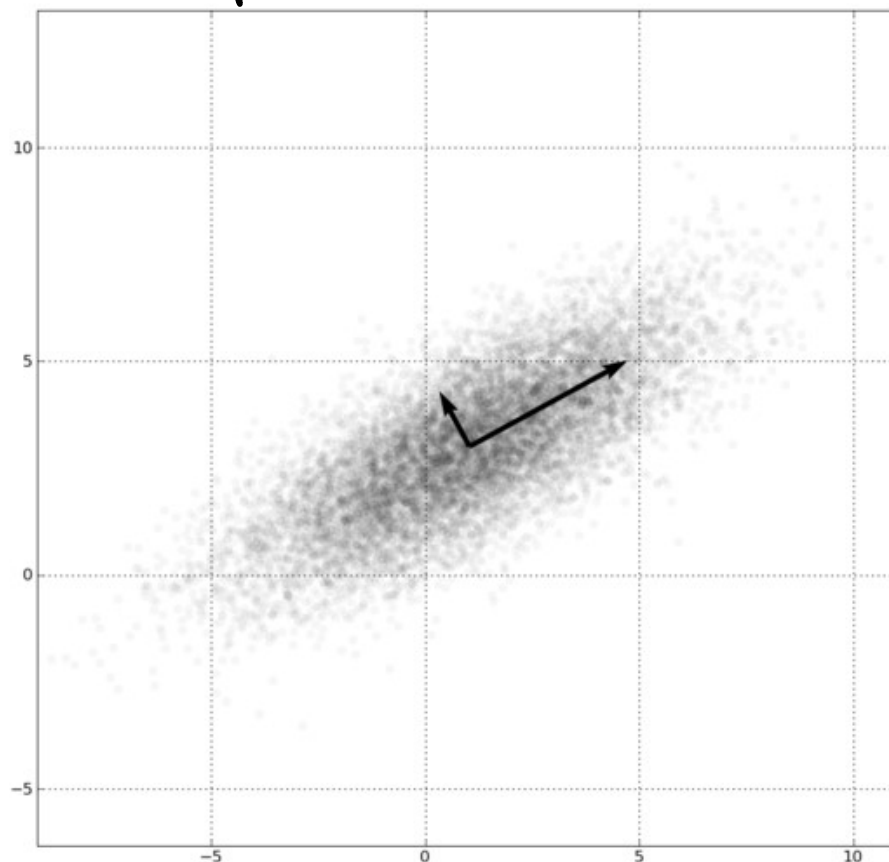
$$\text{Now } P_A x = AA^T x = \underbrace{AA^T A}_{= A} y$$

$$= Ay = x. \quad \text{"A via Moore-Penrose (i)}$$

★ Principal Component Analysis (PCA)

(a.k.a. Karhunen-Loève Transform)
is a data analysis technique that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of linearly uncorrelated variables called "principal components."

2D example (from Wikipedia)



One can understand PCA using SVD! But before doing so, we need a bit of Statistics.

Suppose we are given a set of vectors (observations)

often these \rightarrow are viewed as n realizations of some stochastic process.

x_1, x_2, \dots, x_n and each $x_j \in \mathbb{R}^d$. d : could be huge (ex. a face image database).

Let $X := [x_1 \ x_2 \ \dots \ x_n] \in \mathbb{R}^{d \times n}$

You know the mean (or average) of this data set

$$\bar{x} := \frac{1}{n} \sum_{j=1}^n x_j$$

And define the **centered** data matrix

$$\tilde{X} := [x_1 - \bar{x} \ x_2 - \bar{x} \ \dots \ x_n - \bar{x}]$$

Note: $\tilde{X} = X \left(I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right)$

\hookrightarrow Good exercise!

Now the **sample covariance matrix** S is defined as

$$S := \frac{1}{n} \tilde{X} \tilde{X}^T \in \mathbb{R}^{d \times d}$$

S_{ij} indicates the **covariance** or **mutual correlation** between the i th and j th entries of data vectors.

PCA is nothing but an eigenvalue decomposition of S , i.e.,

$$S = \Phi \Lambda \Phi^T, \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$$

Let's sort λ_i 's as $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$
Because $S^T = S$, and $S = \frac{1}{n} \tilde{X} \tilde{X}^T$,
we can show that $\lambda_i \geq 0, 1 \leq i \leq d$.

$$\Phi = [\Phi_1 \dots \Phi_d] \in \mathbb{R}^{d \times d}$$

is a matrix containing the eigenvectors.
Also thanks to $S^T = S$, Φ is an
orthogonal matrix whose columns
form an ONB of \mathbb{R}^d .

The change of the bases from
 $[e_1 \dots e_d]$ to $[\Phi_1 \dots \Phi_d]$
is achieved simply by $\Phi^T \tilde{X}$.

$\Phi_j^T \tilde{X}$ is called the j th principal
components of X .

PCA was known for a long time,
e.g., since the time of Pearson (1901)
and Hotelling (1933).

Those days, the measurement dimension
 d was much smaller than the number
of samples n , i.e. $d \ll n$.

This is called the "classical" setting.

Ex. 5 exam scores of 2000 students
 $d=5, n=2000$.

Due to the advent of computers and
sensor technology, now we often have
 $d \gg n$, the "neo-classical" setting.

Ex. The face database: $d=128^2, n=143$.